

NS homeshopping shop+ Sales prediction & optimization

5우 좋습니다

Contents Writer:

박준민(조장)
jonah12345@naver.com

김윤환
hellohuman@naver.com

박대한
bighan96@naver.com

정규형
tazan852@naver.com

최종문
chlwhd97@gmail.com

September 2020

Agenda

I. 과제 정의

II. 과제 상세

1. 데이터 전처리
2. 탐색적 자료 분석
3. 외부변수 선정
4. 변수 pool 설정 / Feature Engineering
5. 모델링
6. 최적화 방안

Agenda

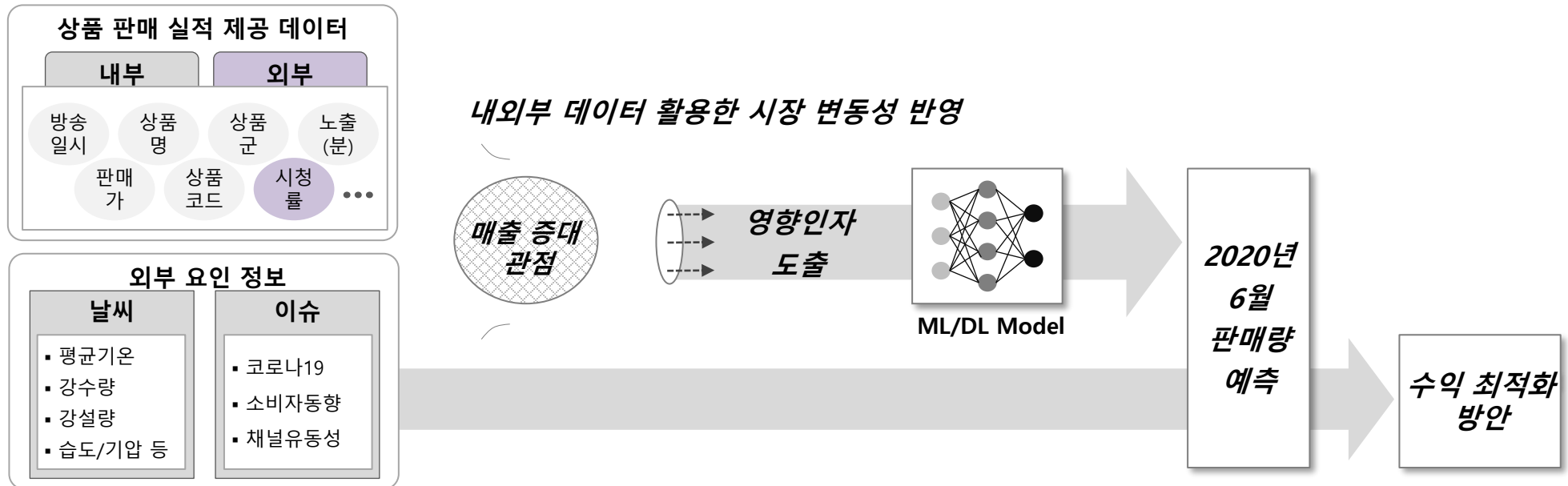
I. 과제 정의

II. 과제 상세

1. 데이터 전처리
2. 탐색적 자료 분석
3. 외부변수 선정
4. 변수 pool 설정 / Feature Engineering
5. 모델링
6. 최적화 방안

I. 과제 정의

- NS Shop+편성데이터(NS홈쇼핑) 를 활용하여 방송편성표에 따른 판매실적을 예측하고, 최적 수익을 고려한 요일별/ 시간대별 / 카테고리별 편성 및 최적화 방안(모형) 제시



추진과제

1 데이터 전처리

- 판매량을 나타내는 target 변수 설정
- 방송일시 카테고리화
- 노출(분) NA imputation

2 모델링

- EDA → 변수 Pool → Engineering 과정 통해 정의된 데이터로 판매량 예측 모델 개발
- model competition을 통한 최종 모델 선정

3 판매량 예측

- 2020년 6월 판매량 예측

4 수익 최적화 방안

- 월별/요일별/시간대별 / 상품군별 최적 편성 제안
- 매출 증대를 위한 효자 상품 제안

I. 과제 정의 – 데이터 개요

2019년 1월 1일부터 2019년 12월 31일 동안 방송된 NS Shop+의 홈쇼핑 편성 정보, 판매 상품 정보, 총 판매액

제공 데이터

NS Shop+ 2019.01.01~12.31								
1	방송일시	노출(분)	마더코드	상품코드	상품명	상품군	판매단가	취급액
2	2019-01-01 6:00	20	100346	201072	테이트 남성 셀린트3중	의류	39,900	2,099,000
3	2019-01-01 6:00		100346	201079	테이트 여성 셀린트3중	의류	39,900	4,371,000
4	2019-01-01 6:20	20	100346	201072	테이트 남성 셀린트3중	의류	39,900	3,262,000
5	2019-01-01 6:20		100346	201079	테이트 여성 셀린트3중	의류	39,900	6,955,000
6	2019-01-01 6:40	20	100346	201072	테이트 남성 셀린트3중	의류	39,900	6,672,000
7	2019-01-01 6:40		100346	201079	테이트 여성 셀린트3중	의류	39,900	9,337,000
8	2019-01-01 7:00	20	100305	200974	오모데 레이스 파운데이션 브라	속옷	59,000	6,819,000
9	2019-01-01 7:20	20	100305	200974	오모데 레이스 파운데이션 브라	속옷	59,000	15,689,000
10	2019-01-01 7:40	20	100305	200974	오모데 레이스 파운데이션 브라	속옷	59,000	25,370,000

평가 데이터

NS Shop+ 2020년 6월 편성								
1	방송일시	노출(분)	마더코드	상품코드	상품명	상품군	판매단가	취급액
2	2020-06-01 6:20	20	100650	201971	잭필드 남성 반팔셔츠 4중	의류	59,800	
3	2020-06-01 6:40	20	100650	201971	잭필드 남성 반팔셔츠 4중	의류	59,800	
4	2020-06-01 7:00	20	100650	201971	잭필드 남성 반팔셔츠 4중	의류	59,800	
5	2020-06-01 7:20	20	100445	202278	쿠미투니카 클 레이시 란쥬셰이퍼&팬티	속옷	69,900	
6	2020-06-01 7:40	20	100445	202278	쿠미투니카 클 레이시 란쥬셰이퍼&팬티	속옷	69,900	
7	2020-06-01 8:00	20	100445	202278	쿠미투니카 클 레이시 란쥬셰이퍼&팬티	속옷	69,900	
8	2020-06-01 8:20	20	100381	201247	바비리스 퍼펙트 볼륨스타일러	이미용	59,000	
9	2020-06-01 8:40	20	100381	201247	바비리스 퍼펙트 볼륨스타일러	이미용	59,000	
10	2020-06-01 9:00	20	100381	201247	바비리스 퍼펙트 볼륨스타일러	이미용	59,000	

데이터 설명

- 기간: '19.01 ~ '19.12
- 총 38309개
- 방송일시
- 노출(분) : 연속 방영 시간
- 마더코드 : 상품 대분류
- 상품코드
- 상품명
- 상품군 : 의류, 속옷, 주방, 농수축, 이미용, 가전, 생활용품, 건강기능, 잡화, 가구, 침구 총 11개
- 판매단가
- 취급액

- 기간: '20.06
- 총 2892개
- 취급액 예측

Agenda

I. 과제 정의

II. 과제 상세

1. 데이터 전처리
2. 탐색적 자료 분석
3. 외부변수 선정
4. 변수 pool 설정 / Feature Engineering
5. 모델링
6. 최적화 방안

1. 데이터 전처리

세 개의 파이썬 파일로 전처리, 외부변수 병합 및 파생변수 생성을 일괄적으로 처리하여 효율성 제고

preprocessing.py

- make_count() : 판매량 변수 생성
- log_sales_cnt : 종속변수 판매량 정규화
- del_comma() : 숫자형 변수 콤마 제거
- divide_time() : 방송일시 변수 분리
- holiday_dummy() : 공휴일 변수 생성
- month_order() : 월초, 월말 변수 생성
- pd.merge() : 기상정보 외부변수 병합
- fillna(method='ffill') : 결측값 처리

grouping.py

- make_group() : 모델링 기준 그룹 정보 생성
- make_g1() : 그룹별 데이터프레임 분리
- make_up_ind() : unit_price_group 변수 생성
- Make_cpi_csi() : 소비자물가지수 외부변수

spliting.py

- variable_selection() : 모델링 최종변수 선정
- product_name : 그룹1 상품명 변수 가공

제공데이터 df_train

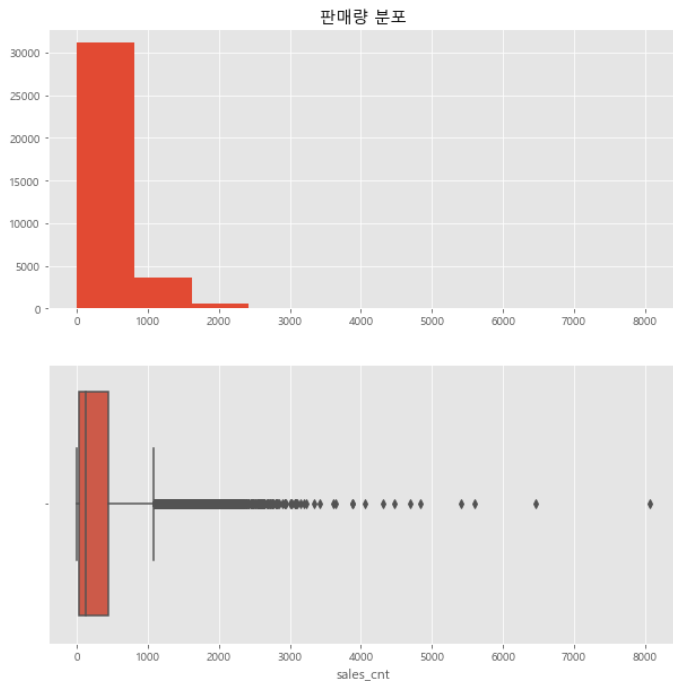
	month	min	grp	cnt
0				
1				
2				
3				

평가데이터 df_test

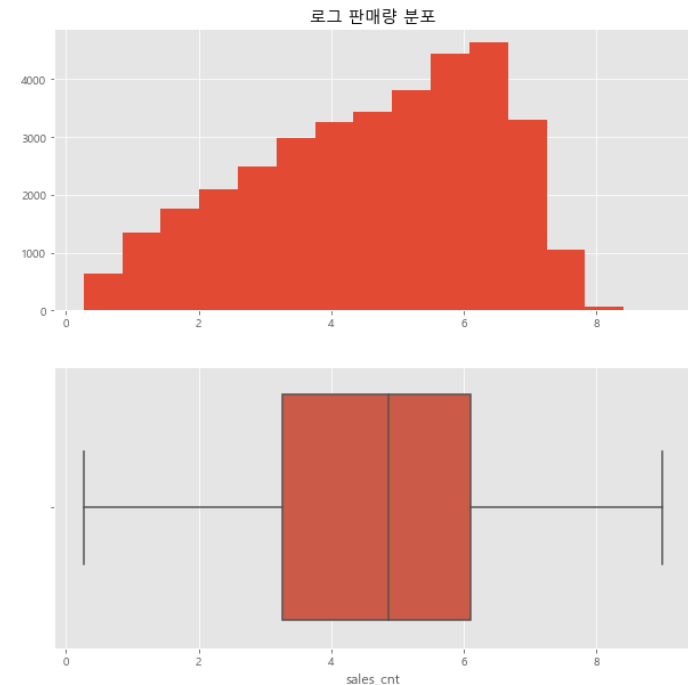
	month	min	grp	cnt
0				
1				
2				
3				

<Target 변수> 판매량 (sales_cnt) = 판매액(sell_price) 을 단위가격(unit_price)으로 나눈 값

원본 판매량



로그변환 판매량

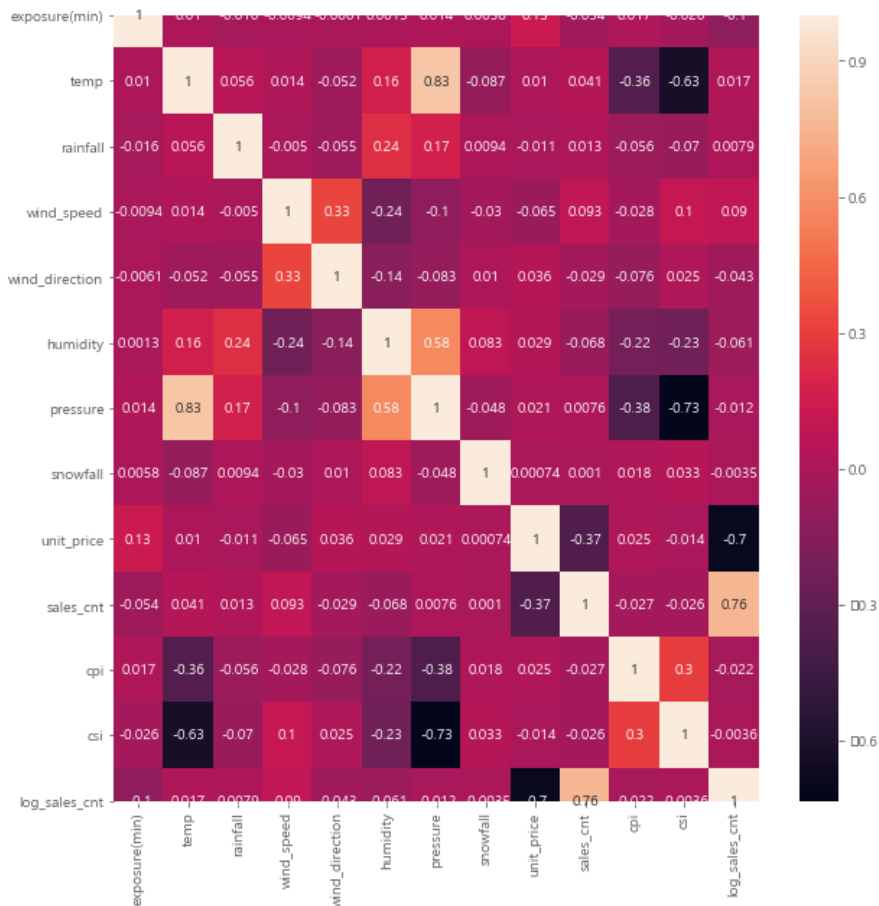


변환 전 판매량은 지나치게 밀집되어 있어 로그변환한 판매량을 쓰기로 한다.

2. 탐색적 자료 분석 - 판매량

피어슨 상관계수 그래프를 통한 변수들간 상관관계 파악

피어슨 상관계수 그래프

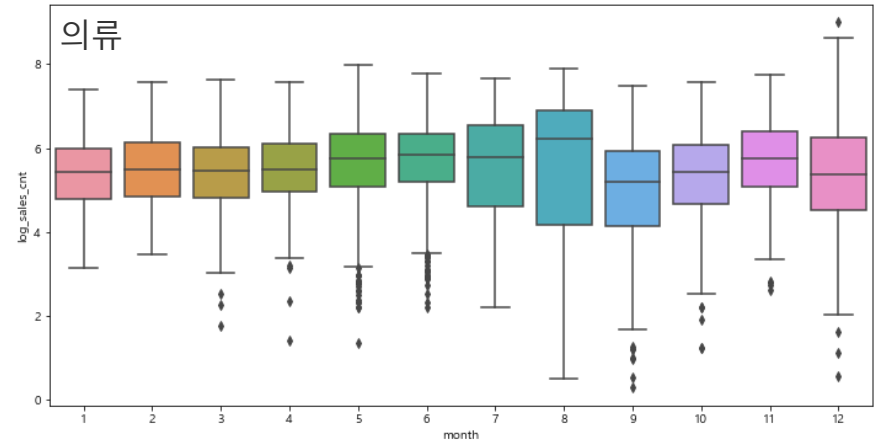
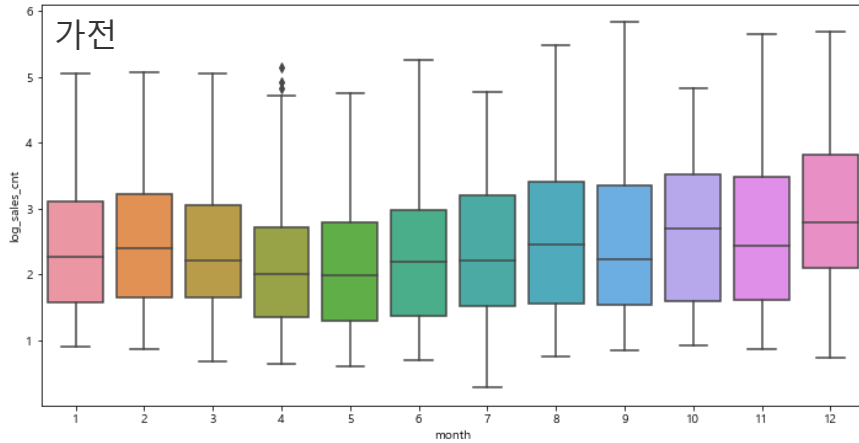


분석 결과

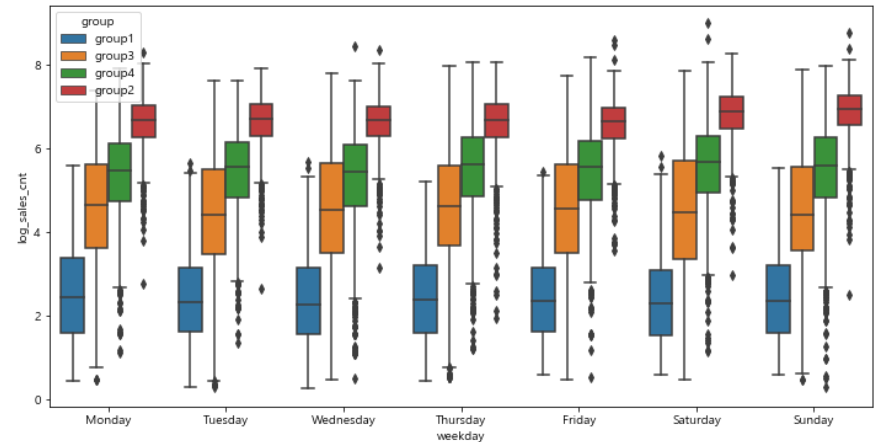
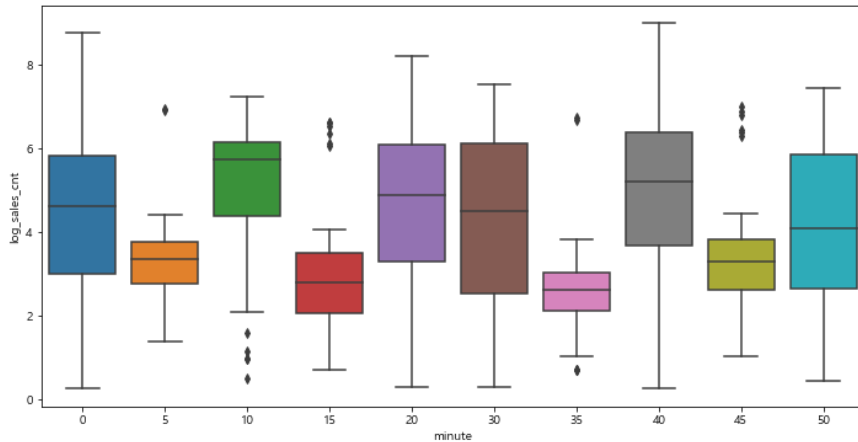
- 상품의 단가와 판매량이 가장 큰 음의 상관관계를 갖는다 (-0.37)
(단가가 높을수록 적게 팔린다)
- 판매량을 로그변환했을 때
단가와 판매량의 상관계수가 더 커진다
(-0.37 -> -0.7)
- 단가를 제외한 변수들은
판매량과의 상관계수가 매우
작은 것으로 보아 선형관계가
없는 것으로 보인다

2. 탐색적 자료 분석 - 시간 변수

월, 분, 요일 등의 시간 변수를 종속변수에 대해 EDA 진행 및 유의미한 변수 추출



- 월별, 상품군별 조합에 따라 판매량 추이가 달라짐

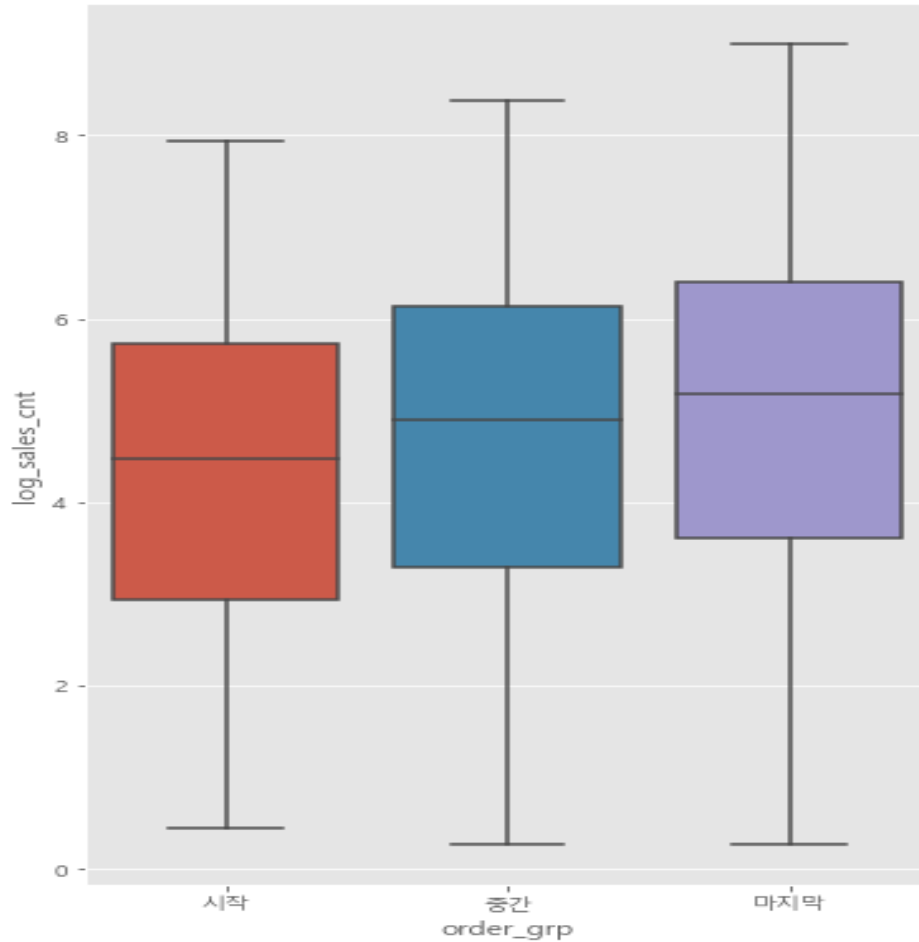


- 그룹별, 요일별에 따른 판매량 추이도 명확히 두드러짐

2. 탐색적 자료 분석 - 방송 순서 그룹화

같은 상품의 방송 송출 순서에 따른 로그-판매량의 분포

방송 순서 그룹(order_group) 별 로그-판매량(log sales_cnt)



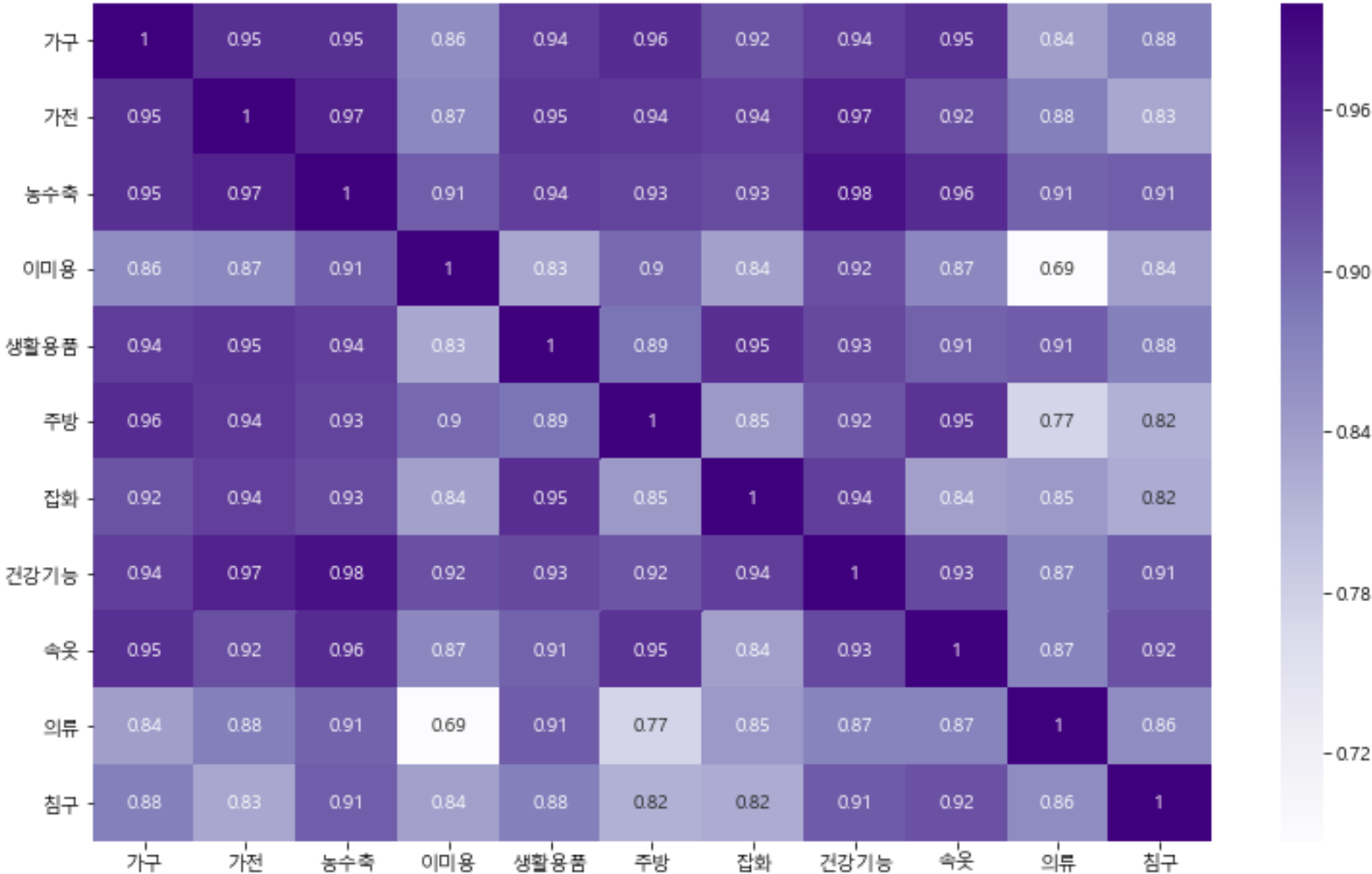
그룹화 결과

- 같은 상품을 계속 판매 하는 경우 마지막 방송에서 판매량이 증가
- 같은 상품의 연속 방송의 경우 방송 순서에 따라 시작, 중간, 마지막으로 표시
- 최종적으로 방송순서에 따른 그룹변수 추가

2. 탐색적 자료 분석 – 상품군 그룹화

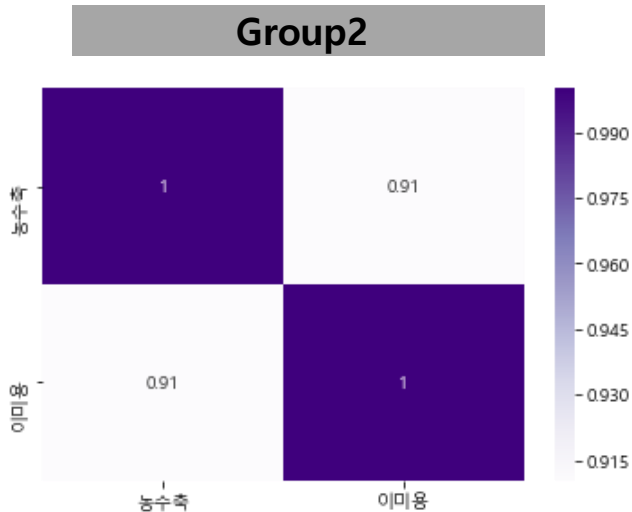
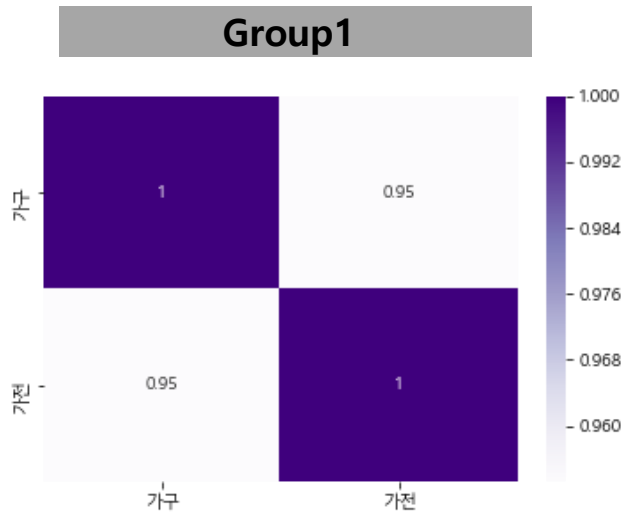
판매량(sales_cnt)에 대한 상품군 상관관계 분석을 통해 그룹화 진행

상품군 별 변수 간 상관관계

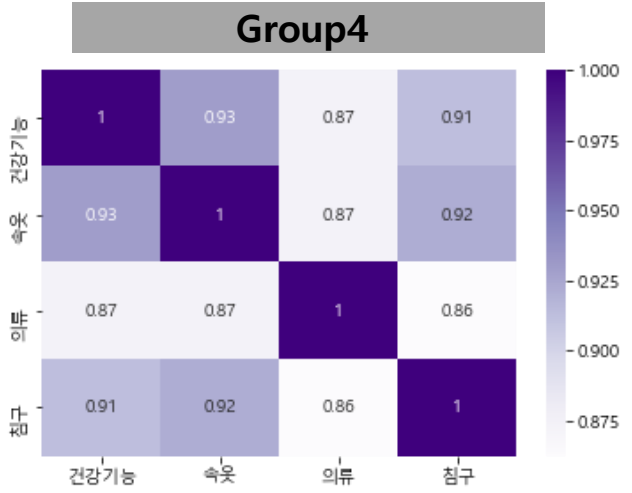
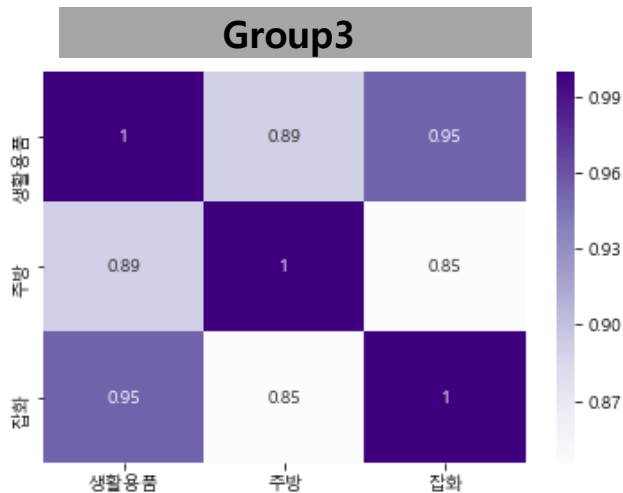


2. 탐색적 자료 분석 - 상품군 그룹화

최종 그룹



그룹명	상품군
group1	가구, 가전
group2	농수축, 이미용
group3	생활용품, 주방, 잡화
group4	건강기능, 속옷, 의류, 침구

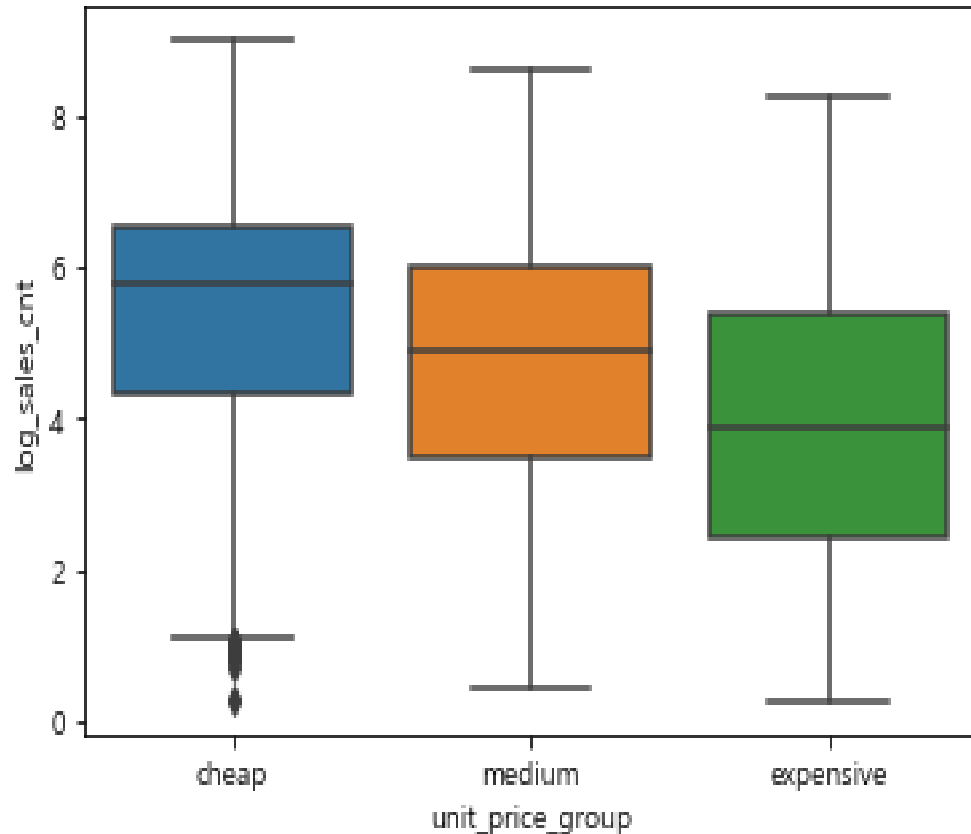


↓
그룹별
모델 구축

2. 탐색적 자료 분석 - 단위 가격 그룹화

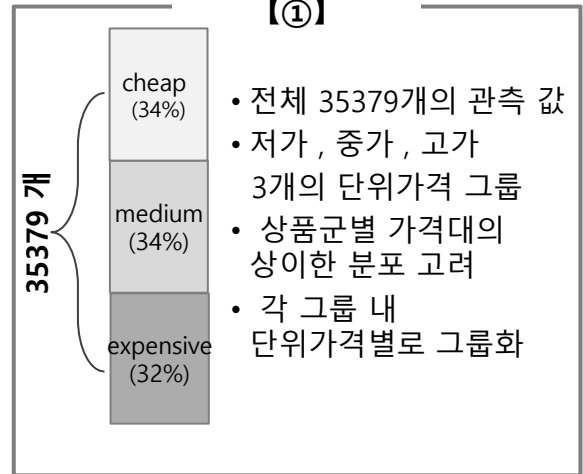
그룹별로 단위가격의 33%, 66% quantile 값을 기준으로 정하여 순서 그룹과 단위 가격별로 저가, 중가, 고가 그룹을 형성

단위가격 그룹(unit_price_group) 별 로그-판매량(log sales_cnt)



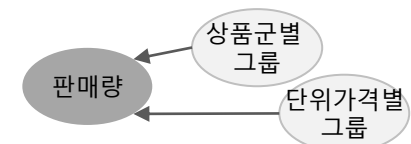
그룹화 결과

【①】



【②】

- 그룹화한 결과 log_sales_cnt별로 상이한 분포를 보여 단위 가격별 그룹화 변수가 타겟변수인 sales_cnt에 유의할 것으로 판단
- 추후 모델링시 최종 변수로 선택



3. 외부변수 설정 - 기상정보 / 이슈정보

기상정보

이슈정보

배경

- 홈쇼핑 시청률은 재택률에 기반한다고 판단
- 재택률 : 집에 머무르는 비율
- 재택률에 영향을 끼치는 기상정보를 외부변수로 설정

- 소비자 동향 및 소비자 심리가 전체 판매액에 영향을 끼칠거라 판단
- 최근 가장 큰 이슈는 코로나 19 관련 소비자 동향이 소비자 관련 지표에 담겨있다고 판단

변수명
및
내용

방송된 시각기준
temp : 기온(°C)
rainfall : 강수량
wind_speed : 풍속(m/s)
wind_direction : 풍향(16방위)
humidity : 습도
pressure : 증기압(hPa)
spot_pressure : 현지기압(hPa)
sea_level_pressure : 해면기압(hPa)
snow_fall : 적설량(cm)

방송된 시각기준
csi : 소비자 동향 지수
cpi : 소비자 물가 지수

정보 출처

- 기상자료개방포털(data.kma.go.kr)
- 서울 기준

- e-나라지표 (www.index.go.kr)

4. 변수 pool 설정 / Feature Engineering

탐색적 자료 분석으로 파생변수를 선정하고 외부변수를 추가하여 최종 변수 pool 선정

가공 데이터

	datetime	year	month	day	hour	minute	weekday	holiday	month_order	order_grp	exposure(min)	mother_cd	product_cd	product_name
0	2019-01-01 06:00:00	2019	1	1	6	0	Tuesday	1	초	시작	20.0	100346	201072	테이트 남성 셀 린니트3종
1	2019-01-01 06:00:00	2019	1	1	6	0	Tuesday	1	초	시작	20.0	100346	201079	테이트 여성 셀 린니트3종
2	2019-01-01 06:20:00	2019	1	1	6	20	Tuesday	1	초	중간	20.0	100346	201072	테이트 남성 셀 린니트3종

	product_grp	temp	rainfall	wind_speed	wind_direction	humidity	pressure	spot_pressure	sea_level_pressure	snowfall
0	의류	-7.9	0.0	1.3	290.0	60.0	2.0	1023.6	1034.9	0.0
1	의류	-7.9	0.0	1.3	290.0	60.0	2.0	1023.6	1034.9	0.0
2	의류	-7.9	0.0	1.3	290.0	60.0	2.0	1023.6	1034.9	0.0

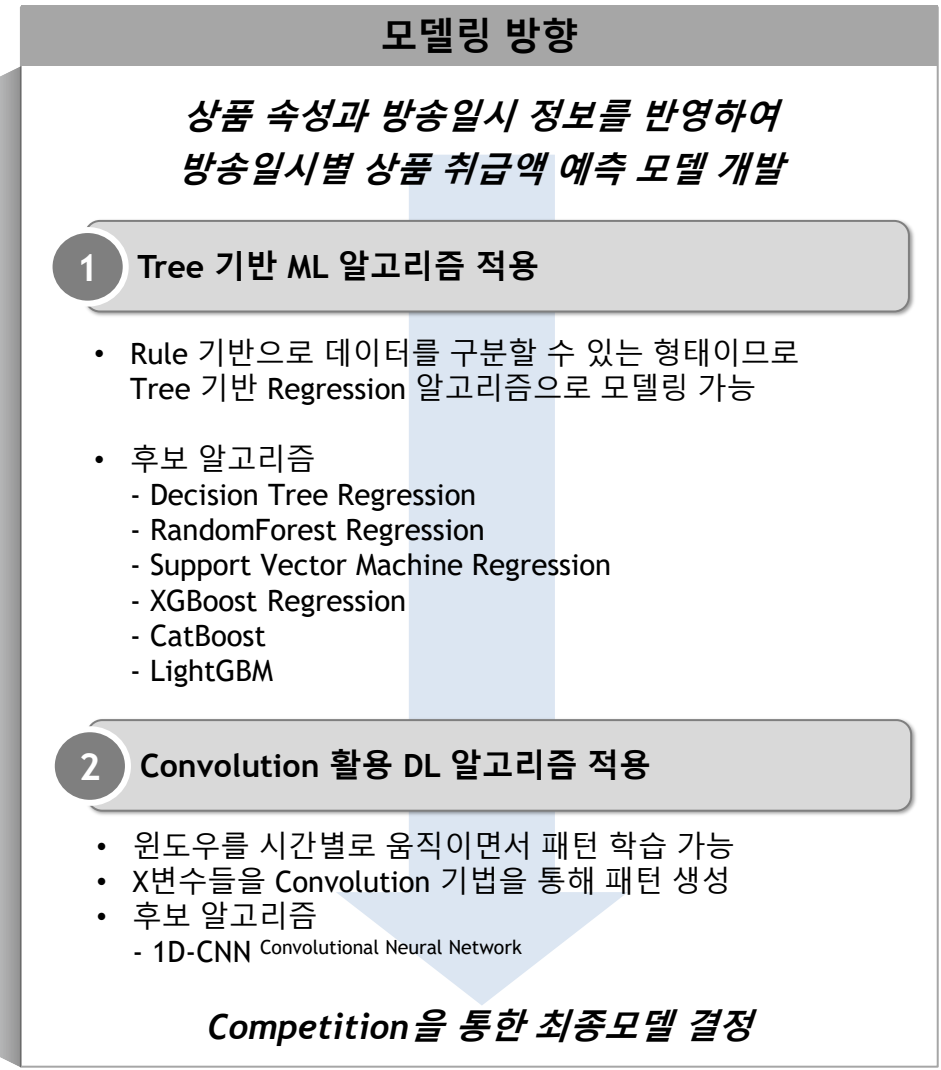
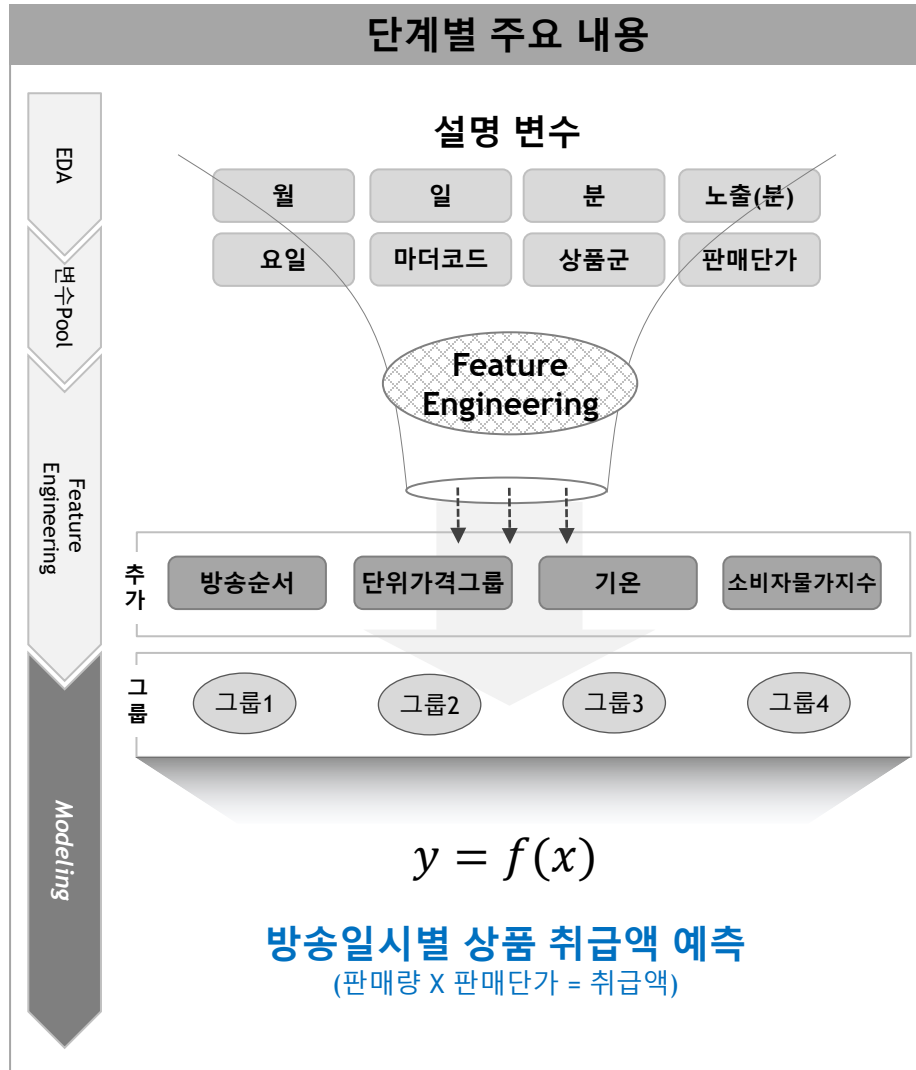
	unit_price	sell_price	sales_cnt	group	unit_price_group	cpi	csi
0	39900.0	2099000.0	52.606516	group4	cheap	100.8	97.5
1	39900.0	4371000.0	109.548872	group4	cheap	100.8	97.5
2	39900.0	3262000.0	81.754386	group4	cheap	100.8	97.5

변수 설명

- 기간: '19.01 ~ '19.12
- 판매가 0인 행, 취급액 없는 행 삭제 등 총 38309개
- dateline을 year, month, day, hour, minute으로 나눔
- weekday : dateline을 활용한 요일
- holiday : 휴일 여부 (휴일이면 1, 아니면 0)
- month_order : 월초, 중순, 월말
- order_group : 같은 상품을 나눠서 방송 했을 때, 시작, 중간, 마지막
- 기온, 강수량, 풍속, 풍향, 습도, 기압, 강설량 등의 기상정보
- sales_cnt : 취급액/단가
- group : 상품군을 4개로 나눔
- unit_price_group : group 내 가격(cheap, medium, expensive)
- cpi : 소비자 물가지수
- csi : 소비자 동향지수

5. 모델링

판매량에 영향을 미치는 유의미한 변수를 EDA 과정을 통해 확인하여, 후보변수 중 16개 변수를 변수Pool 대상으로 정의함. 모델링은 그룹을 기준으로 Competition을 통해 최종모델 결정.



5. 모델링

Competition을 통해 그룹별 최적 모델을 선정하고 Bayesian Optimization을 통해 Hyperparameter Tuning 진행

1

MAPE 기준 Best Model 선정 : CatBoost, LightGBM

	DecisionTree	RandomForest	XGBoost	CatBoost	LighGBM	1D-CNN
그룹1	54.32	52.30	57.49	<u>50.45</u>	<u>48.46</u>	59.39
그룹2	33.48	29.80	57.49	<u>29.58</u>	<u>25.98</u>	38.12
그룹3	53.90	<u>42.97</u>	47.14	48.35	<u>39.25</u>	52.32
그룹4	55.39	49.03	62.22	<u>50.70</u>	<u>48.73</u>	55.35

2

Best Model에 Bayesian Optimizatoin 적용

```

bayesmodel2 = BayesianOptimization(f = model2_mape, pbounds = pbounds2, verbose = 2, random_state = 4)
bayesmodel2.maximize(init_points=2, n_iter = 100)

```

iter	target	baggin...	learni...	max_depth	min_ch...	n_esti...	num_le...	subsam...
[LightGBM] [Warning] boosting is set=gbd								
[LightGBM] [Warning] bagging_freq is set=8,								
1	-27.05	8.703	0.1047	48.66	21.73	1.547e+0	11.59	0.9763
[LightGBM] [Warning] boosting is set=gbd								
[LightGBM] [Warning] bagging_freq is set=0,								
2	-29.48	0.05607	0.0753	22.3	23.6	796.5	43.29	0.9834
[LightGBM] [Warning] boosting is set=gbd								
[LightGBM] [Warning] bagging_freq is set=4,								
3	-28.35	4.634	0.1246	46.22	15.09	1.55e+03	13.04	0.6569
[LightGBM] [Warning] boosting is set=gbd								
[LightGBM] [Warning] bagging_freq is set=8,								
4	-29.2	8.403	0.09335	11.84	20.91	673.5	35.79	0.6187
[LightGBM] [Warning] boosting is set=gbd								
[LightGBM] [Warning] bagging_freq is set=5,								
5	-28.45	5.312	0.1492	46.85	20.68	1.544e+0	12.03	0.8344

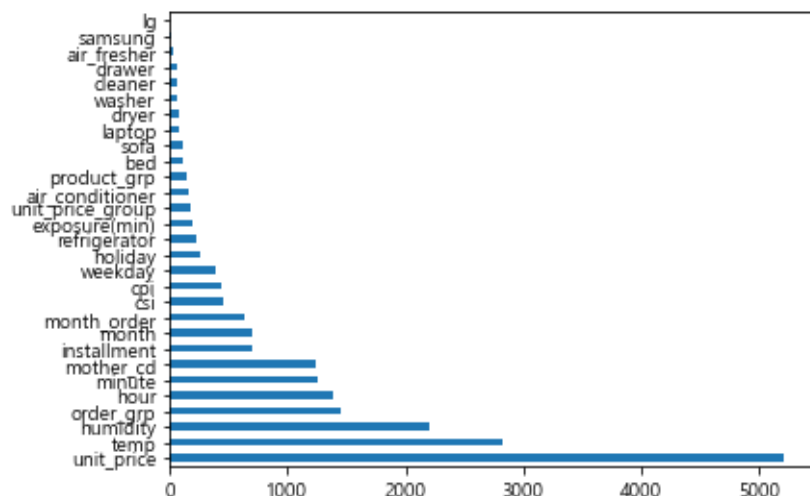
Final Validation Score

	CatBoost	LightGBM
그룹1	49.86	<u>47.60</u>
그룹2	28.01	<u>25.05</u>
그룹3	39.18	<u>37.23</u>
그룹4	49.31	<u>45.40</u>

최종 모델로 LightGBM 선정

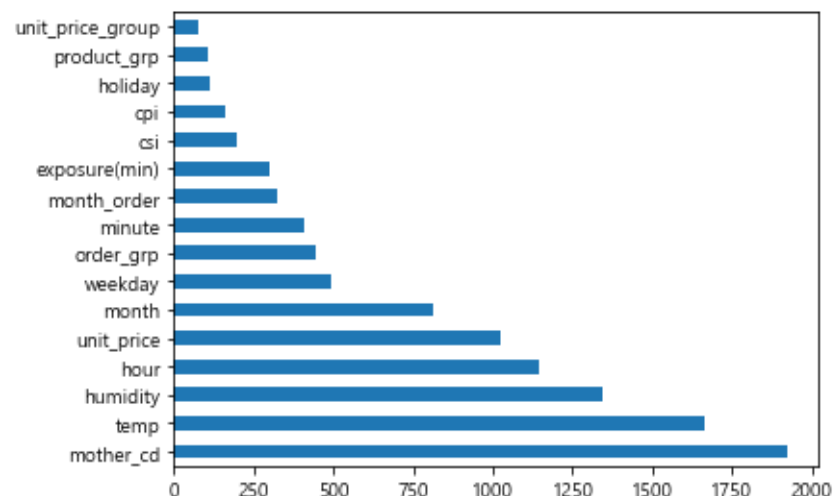
그룹 별 모델의 Feature Importance

Feature Importance (Group1)



- 첫 번째 그룹의 모델은 상품의 단가, 기온, 습도, order_group 순으로 영향을 많이 받음

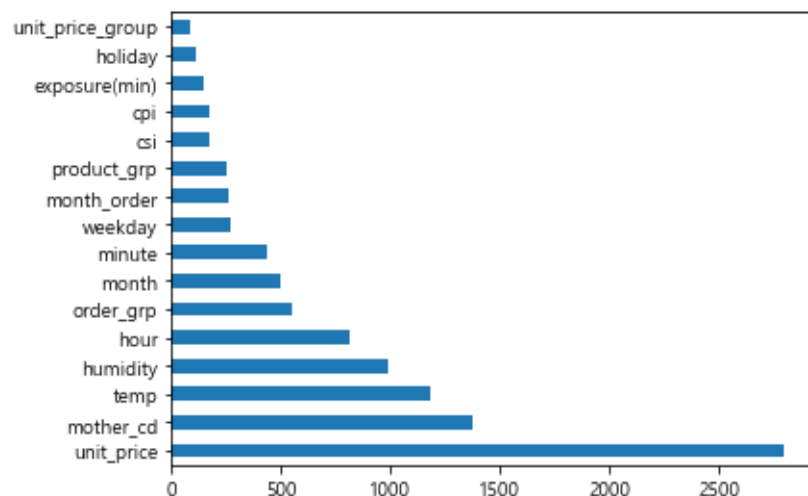
Feature Importance (Group2)



- 두 번째 그룹의 모델은 상품의 mother_cd, 기온, 습도, 시간, 상품의 단가 순으로 영향을 많이 받음

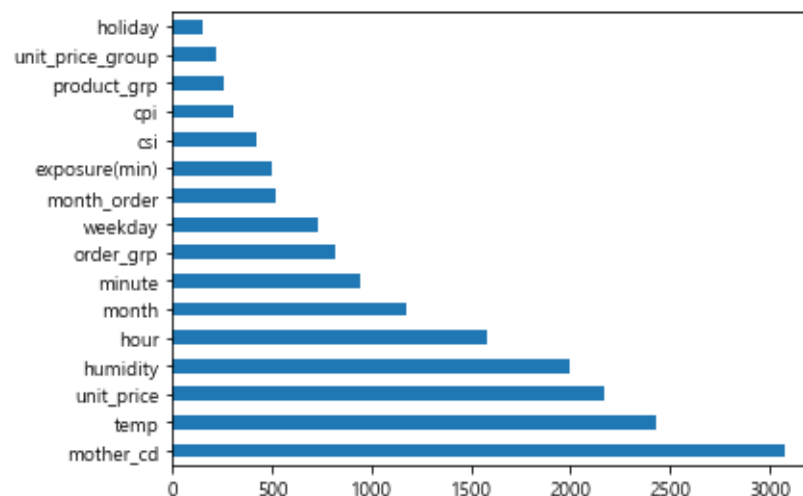
그룹 별 모델의 Feature Importance

Feature Importance (Group3)



- 세 번째 그룹의 모델은 상품의 단가, mother_cd, 기온, 습도, 시간 order_group 순으로 영향을 많이 받음

Feature Importance (Group4)

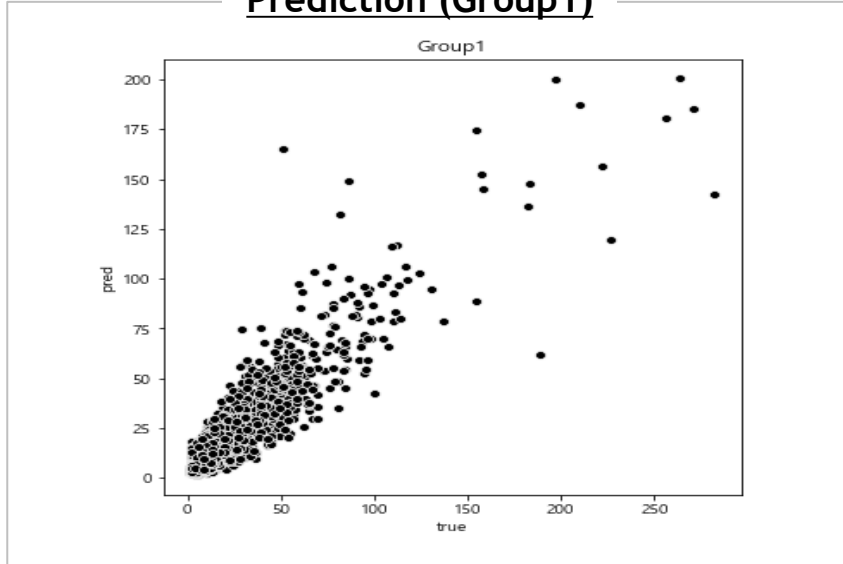


- 네 번째 그룹의 모델은 상품의 mother_cd, 기온, 상품의 단가, 습도, 시간 순으로 영향을 많이 받음

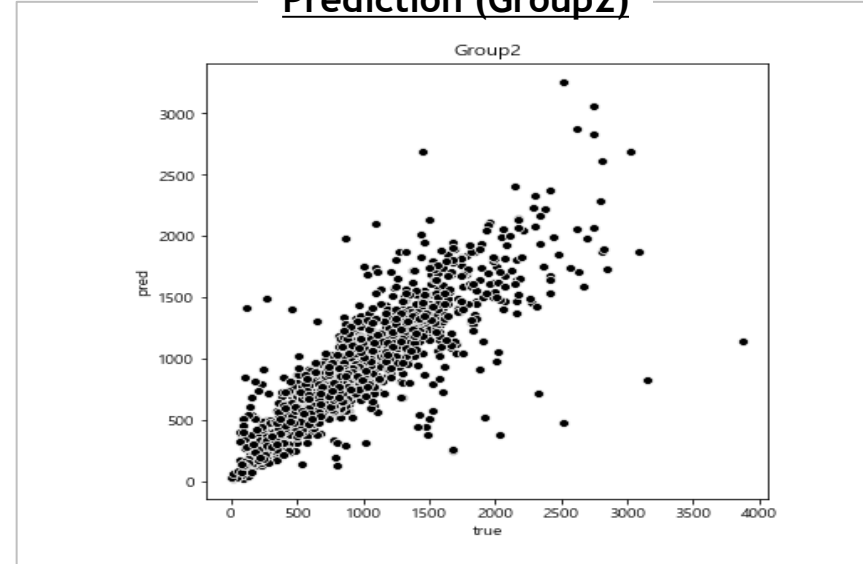
5. 모델링 – 최종 모델 lightGBM

train_test_split 한 후 test set에서의 예측
Scatter-Plot (X축은 true-value, Y축은 prediction-value)

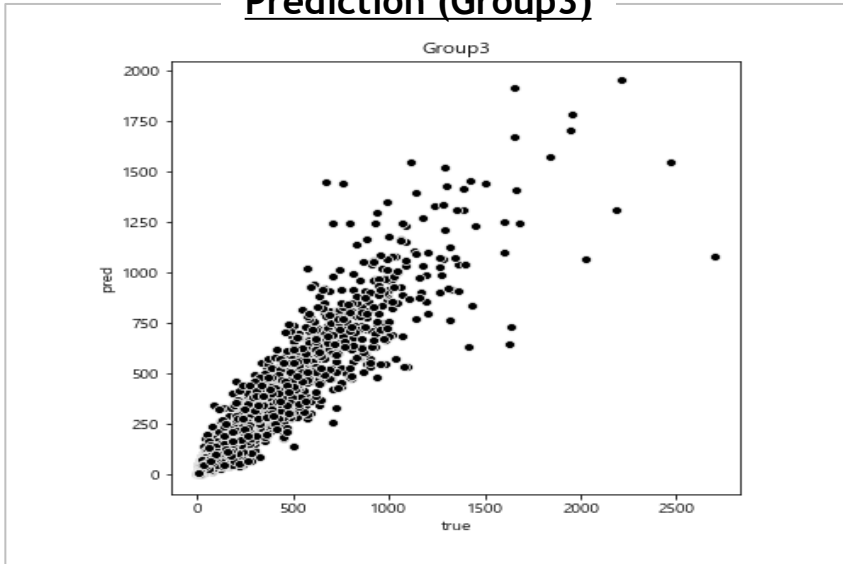
Prediction (Group1)



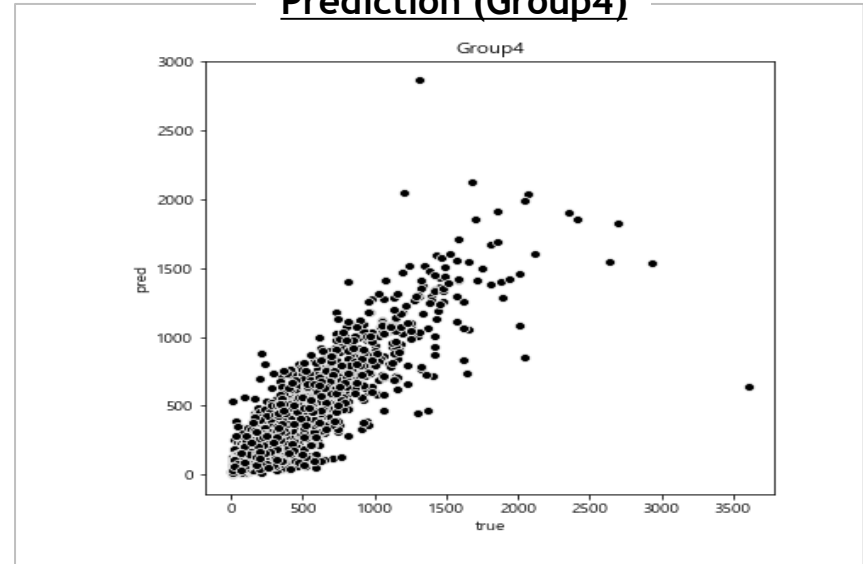
Prediction (Group2)



Prediction (Group3)

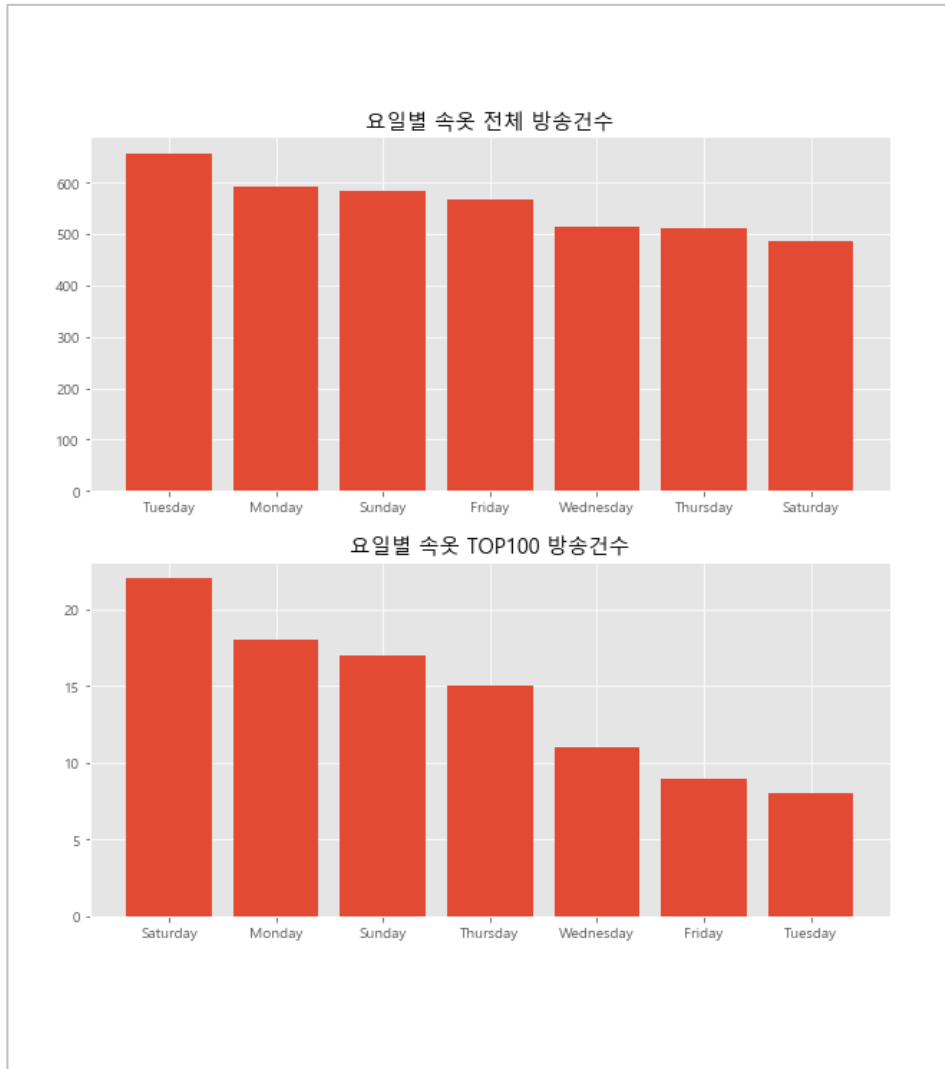


Prediction (Group4)

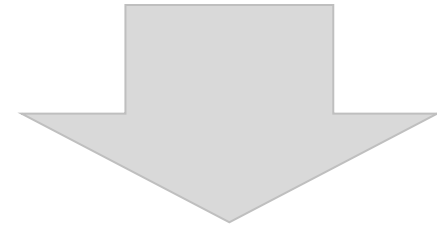


6. 최적화 방안 - 상품군별 최적화 <속옷>

요일 별 속옷의 전체 방송건수와 판매량 기준 TOP100개의 요일 별 속옷의 방송건수 비교

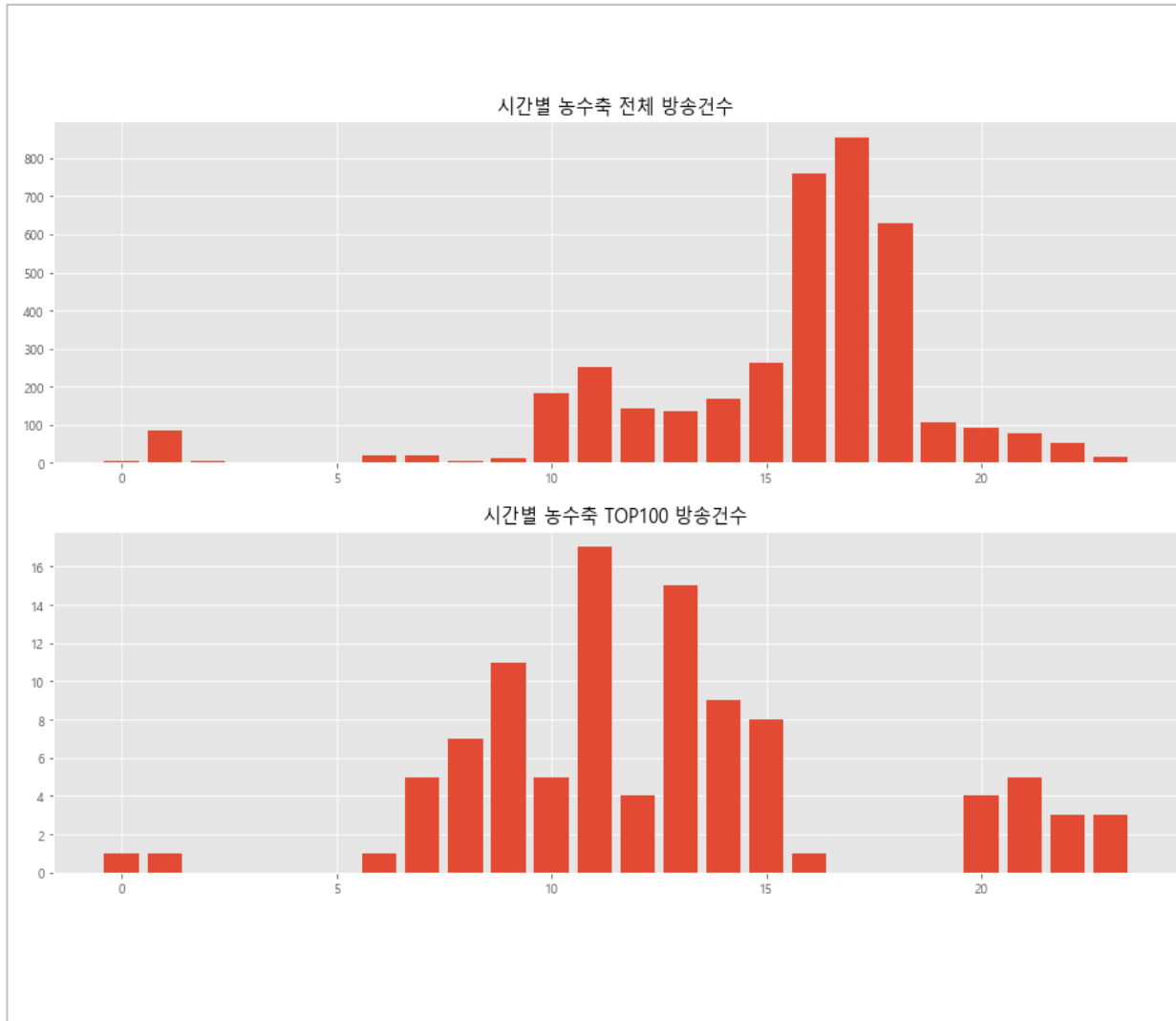


- 전체 속옷의 토요일 판매건수는 다른 요일에 비해 비중이 적다. 하지만 판매량 기준 TOP100의 속옷은 토요일에 방송된 비율이 높았다.
- 반대로 전체 속옷의 화요일 판매건수는 다른 요일에 비해 비중이 높다. 하지만 판매량 기준 TOP100의 속옷은 화요일에 방송된 비율이 제일 낮았다.

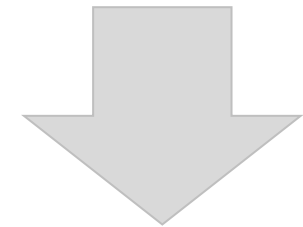


- 주력 속옷 상품 혹은 인기 있는 속옷 상품을 토요일에 더 배치하면 매출이 증가 할 것이다.
- 반대로 주력 상품 혹은 인기 있는 속옷 상품을 화요일을 피해서 배치해야 더 높은 매출을 기록할 수 있다.

요일 별 농수축의 전체 방송건수와 판매량 기준 TOP100개의 요일 별 농수축의 방송건수 비교



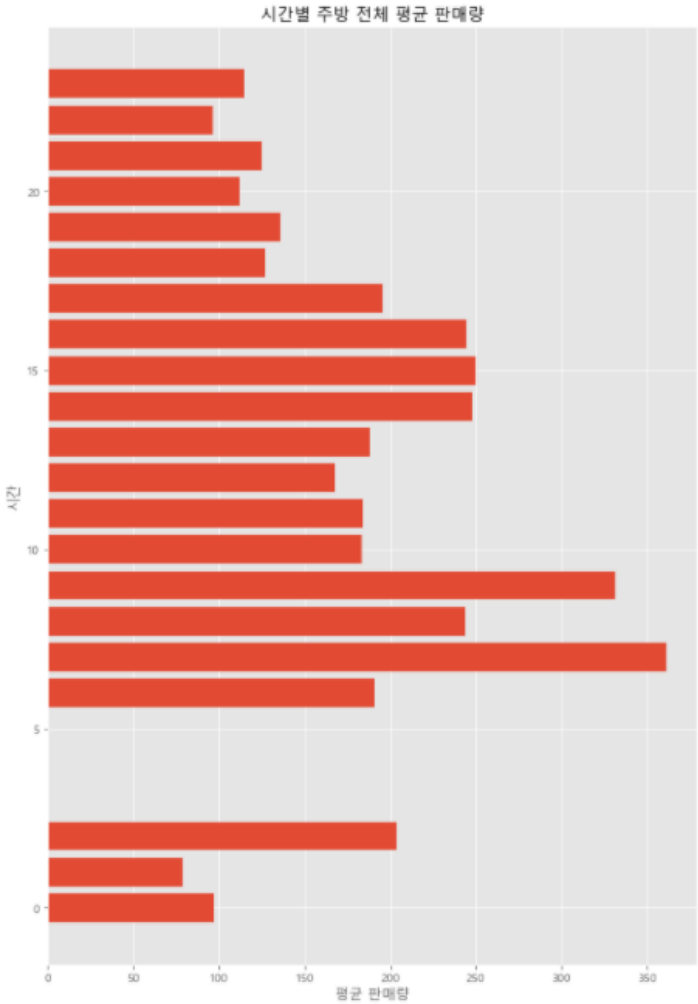
- 전체 농수축 상품의 16시 ~ 18시 방송건수는 다른 시간에 비해 비중이 매우 크다.
- 하지만 판매량 기준 TOP100의 농수축 상품의 16시 ~ 18시 방송건수는 굉장히 적거나 없었다.



- 농수축 상품의 경우 16시 ~ 18시에 배치를 최소화하면 매출이 더욱 증가할 것이다.

6. 최적화 방안 - 상품군별 최적화 <주방>

주방용품의 시간 별 평균 판매량



- 7시 ~ 9시, 14시 ~ 16시에 판매한 주방용품의 평균 판매량이 높다.
-
- 주력 주방용품 혹은 인기 있는 주방용품을 7시 ~ 9시 또는 14시 ~ 16시에 배치하면 매출이 증가 할 것이다.

가전제품의 가격대 별 평균 판매량과 평균 판매액

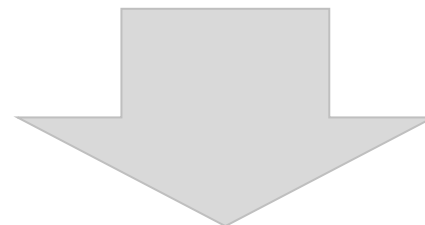
가전제품의 가격대별 평균 판매량

sales_cnt	
unit_price_group	
cheap	54.257131
expensive	9.397112
medium	13.934400

가전제품의 가격대별 평균 판매액

sell_price	
unit_price_group	
cheap	24978325.18
expensive	17828543.60
medium	16638421.52

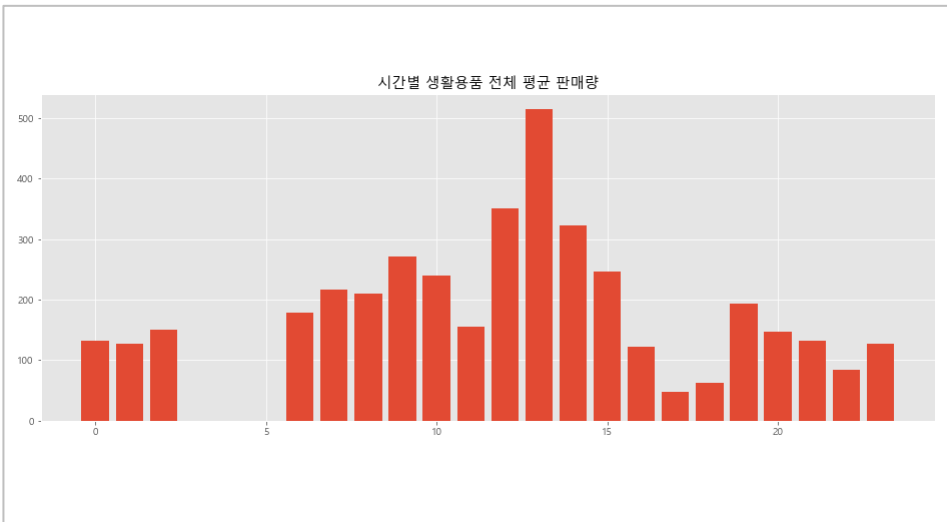
- 저렴한 가격대의 가전제품이 비싸거나 중간 가격대의 가전제품에 비해 평균적으로 더 많이 팔리고 판매액 또한 평균적으로 더 높다



- 저렴한 가격대의 가전제품에 집중해서 판매하면 매출이 더욱 증가할 것이다.

6. 최적화 방안 - 상품군별 최적화 <생활용품>

생활용품의 시간 별 평균 판매량 / 가격대 별 평균 판매량과 평균 판매액



- 12시 ~ 14시의 평균 판매량이 다른 시간대에 비해 높다.



- 생활용품 주력 상품을 12시 ~ 14시에 배치할 경우 더 높은 매출을 기록할 것이다.

가격대 별 평균 판매량
sales_cnt

unit_price_group	
cheap	338.637110
expensive	91.164592
medium	104.565233

가격대 별 평균 판매액
sell_price

unit_price_group	
cheap	19027039.46
expensive	21818785.07
medium	14221059.93

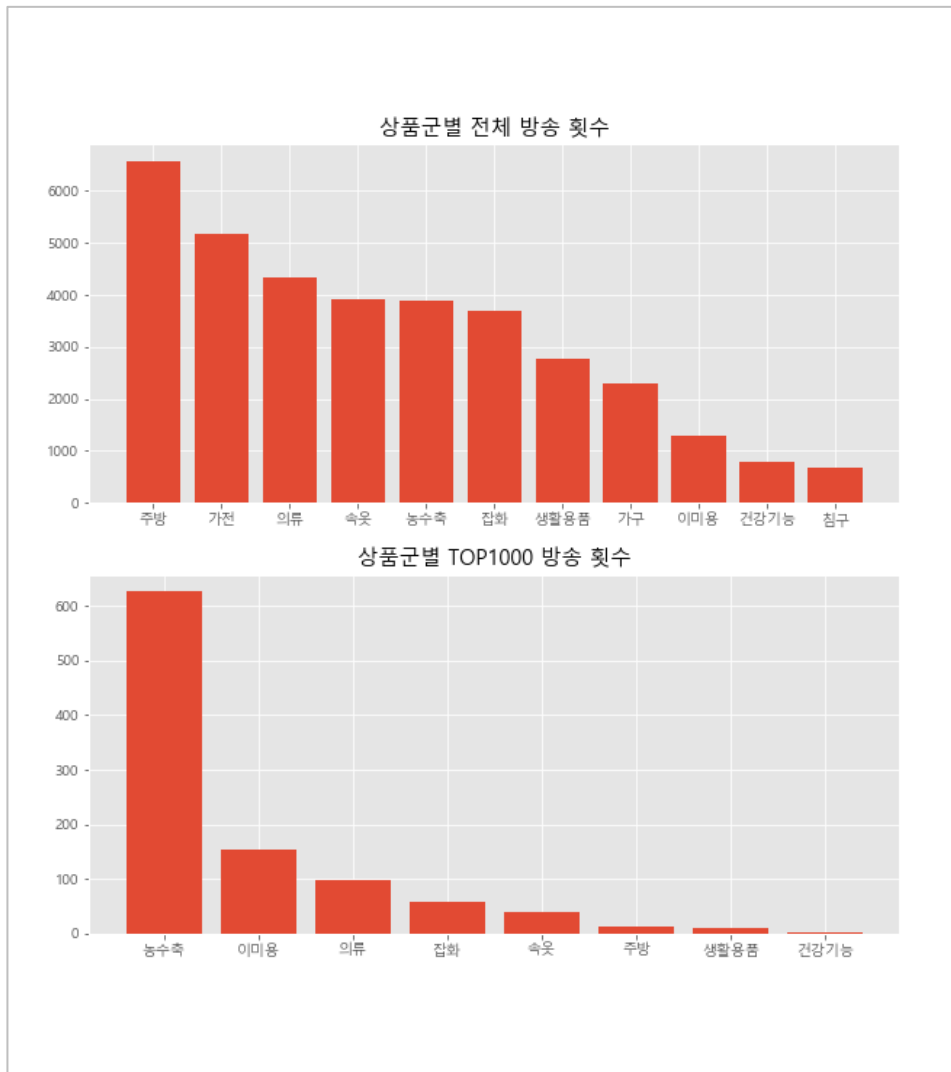
- 저렴한 가격대의 평균 판매량이 가장 높고 평균 판매액은 비싼 가격대가 가장 높다.



- 생활용품의 경우 저가상품과 고급상품을 주력으로 하면 더 높은 매출을 기록할 수 있다.

6. 최적화 방안 – TOP1000 분석을 통한 효자 상품 도출

판매량 기준 상위 1000개의 상품과 그렇지 않은 상품 비교 분석



- 전체 평균 판매량 : 314.797
- TOP1000 상품 평균 판매량 : 1865.923

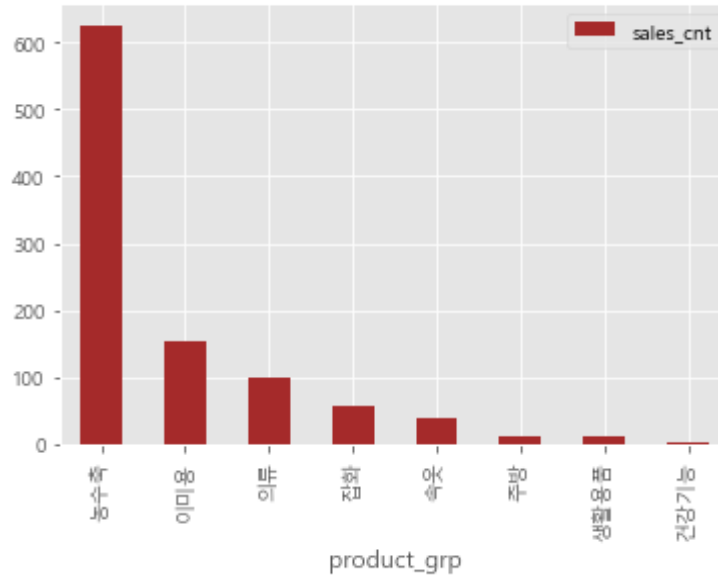
- 전체 방송 횟수는 농수축이 11개의 상품군 중 상위 5번째이지만, 판매량 상위 1000개중에서는 62%로 가장 많은 비중을 차지



- 농수축 상품을 주력으로 판매하여 농수축 상품 대표 홈쇼핑이라는 이미지 각인으로 매출 증대 예상
- 추후 다른 상품군에 대해 공격적인 마케팅을 통해 사업 다각화 가능

6. 최적화 방안 – TOP1000 분석을 통한 효자 상품 도출

TOP1000 기준 상품군 비중



효자상품 상위 27 품목

	mother_cd	product_grp	product_name	rate
16	100699	농수축	고창 꿀 고구마 10kg	1.00
12	100548	농수축	완도꼬마활전복 1.3kg	0.95
15	100637	농수축	영산포숙성홍어회7팩	0.80
0	100435	농수축	우리바다 손질왕고막 24팩	0.53
1	100253	농수축	안동간고등어 20팩	0.46
13	100161	잡화	시스마르스 플렉시블 웨지 폼프스	0.43
4	100698	농수축	강원도양구 간편시래기 + 시래기 들깨 무침	0.42
6	100818	의류	보코 에스닉 앙상블	0.42
5	100450	생활용품	코튼데이 유기농 순면 마스크 KF94 60매	0.38
7	100777	속옷	뷰티플렉스 풍기인건 원피스 2종(8월)	0.36
14	100512	농수축	국내산참조기12팩	0.33
9	100019	잡화	AAD 도트펀칭 컴포트 슬립온	0.33
11	100754	잡화	아가타 소가죽 데이쿠션토퍼	0.33
8	100197	잡화	오델로 여성 겨울모자 3종	0.33
3	100797	농수축	장보고 완도매생이 30개	0.33
17	100236	주방	벨라홀 스마트 멀티포트 1+1 세트	0.33
2	100026	농수축	궁중 손질새우 200미 + 동태포 400g	0.32
10	100829	농수축	참바다손질낙지100미+양념장+연포탕육수	0.30

내용

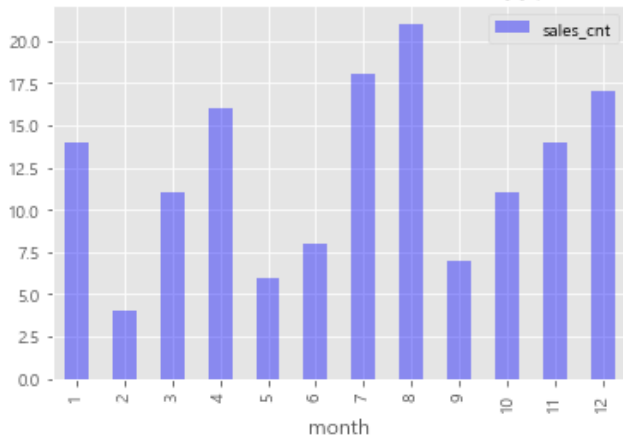
- NS 홈쇼핑 shop + 주력상품답게 농수축이 압도적
- TOP1000에 포함된 상품군별로 효자 상품 선정
- 마더 코드별 총 방송 노출 횟수 대비 TOP1000 포함 비율을 기준으로 0.3이상의 마더코드 호출
- 총 27개의 효자상품 도출

6. 최적화 방안 – 효자 상품 방송 배치 <안동간고등어>

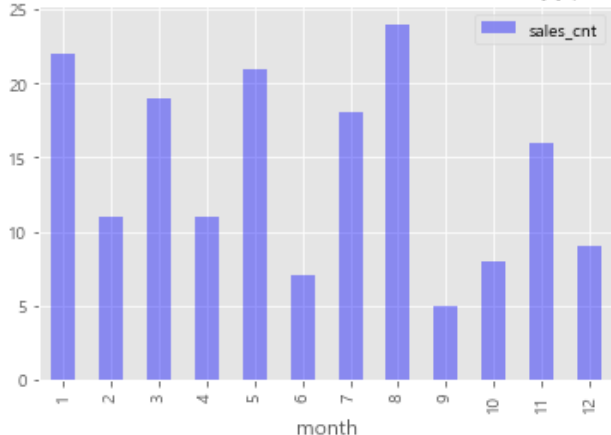
농수축 상품군의 안동간고등어는 제공데이터 기준 318번 방송되었고, 그 중 147번의 방송이 TOP1000에 포함되어 46.2%의 비중을 보임

방송 편성 제안

안동간고등어 TOP1000 월별 판매 횟수



안동간고등어 中 TOP1000 제외 월별 판매 횟수

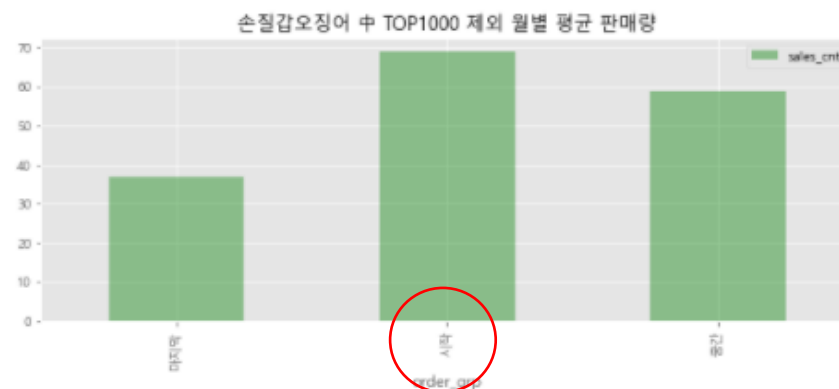
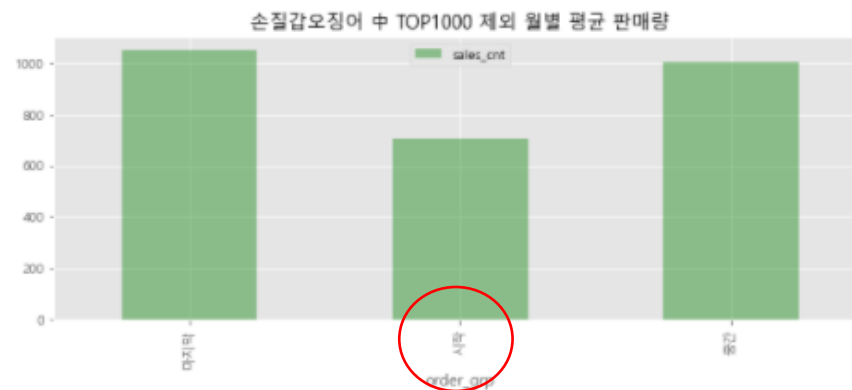


안동간고등어 TOP1000 평균 판매량



- 안동간고등어는 전체 평균 판매량인 314.797보다 월등한 판매량을 보임
- 현재 8월에 많은 판매 방송을 시도하지만 TOP 1000의 월별 판매 횟수와 TOP1000 제외 월별 판매 횟수를 살펴보면 11월, 12월에 판매량을 늘리는 것이 우수한 판매량을 보일 가능성이 큼

방송 편성 제안



내용

- 국내산 손질갑오징어는 방송 노출 기준 처음 방송시에는 TOP1000에 포함되지 않음
- 방송 초반 공략이 새로운 매출 증대의 기회가 될 수 있음
- 초반 할인 이벤트, 연예인 출연 등의 마케팅 전략으로 방송 초반 매출 증대를 기대



End of Document